# Reinforcement learning: An introduction
## Chapter 7

RL Reading Group

# Recap: MC & TD

- ▶ MC: Must wait until the end of the episode to determine the increment to V (St)

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ G_t - V(S_t) \right]$$

- ▶ TD: Form a target and make a useful update using the observed at time t + 1

$$V(S_t) \leftarrow V(S_t) + \alpha \left[ \underbrace{R_{t+1} + \gamma V(S_{t+1})}_{\text{TD target } G_{t:t+1}} - V(S_t) \right]$$
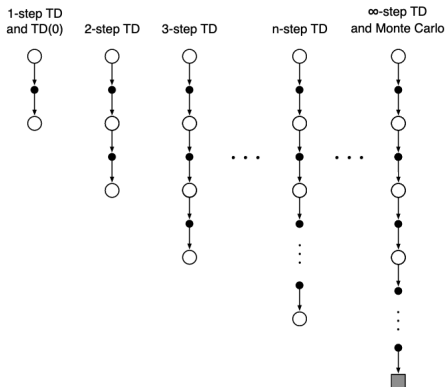
- ▶ NEITHER ARE ALWAYS THE BEST!

# Chapter 7: N-step TD

▶ target for a two-step update is the two-step return

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

▶ target for a n-step update is the n-step return

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

▶ target for a two-step update is the two-step return

$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{t+1}(S_{t+2})$$

▶ target for a n-step update is the n-step return

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n V_{t+n-1}(S_{t+n})$$

▶ The update rule is:

$$V_{t+n}(S_t) \leftarrow V_{t+n-1}(S_t) + \alpha \left[ G_{t:t+n} - V_{t+n-1}(S_t) \right]$$

# Chapter 7: N-step TD

---

**$n$-step TD for estimating $V \approx v_\pi$**

Input: a policy $\pi$
Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer $n$
Initialize $V(s)$ arbitrarily, for all $s \in \mathcal{S}$
All store and access operations (for $S_t$ and $R_t$) can take their index mod $n + 1$

Loop for each episode:
    Initialize and store $S_0 \neq$ terminal
    $T \leftarrow \infty$
    Loop for $t = 0, 1, 2, \ldots$ :
    |   If $t < T$, then:
    |     Take an action according to $\pi(\cdot|S_t)$
    |     Observe and store the next reward as $R_{t+1}$ and the next state as $S_{t+1}$
    |     If $S_{t+1}$ is terminal, then $T \leftarrow t + 1$
    |   $\tau \leftarrow t - n + 1$    ($\tau$ is the time whose state's estimate is being updated)
    |   If $\tau \geq 0$:
    |     $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n,T)} \gamma^{i-\tau-1} R_i$
    |     If $\tau + n < T$, then: $G \leftarrow G + \gamma^n V(S_{\tau+n})$           $(G_{\tau:\tau+n})$
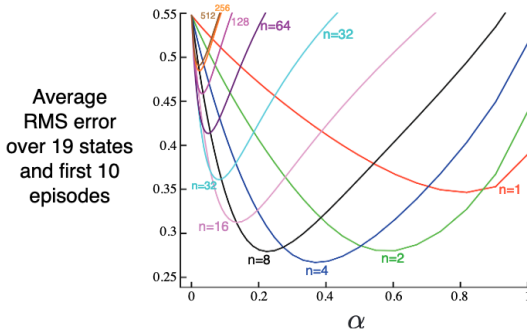    |     $V(S_\tau) \leftarrow V(S_\tau) + \alpha\left[G - V(S_\tau)\right]$
    Until $\tau = T - 1$

# Chapter 7: Random Walk Example Revisited



**Example 6.2   Random Walk**

In this example we empirically compare the prediction abilities of TD(0) and constant-$\alpha$ MC when applied to the following Markov reward process:

start

Average RMS error over 19 states and first 10 episodes

$\alpha$

▶ The target is:

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

▶ The update rule is:

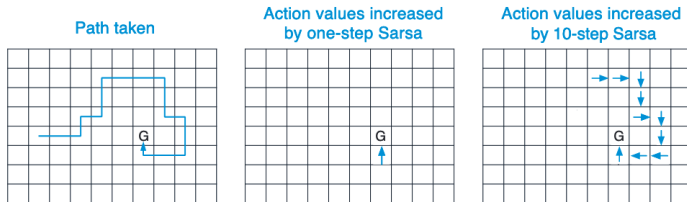$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha \left[ G_{t:t+n} - Q_{t+n-1}(S_t, A_t) \right]$$

# Chapter 7: N-step Sarsa

▶ The target is:

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q_{t+n-1}(S_{t+n}, A_{t+n})$$

▶ The update rule is:

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha \left[ G_{t:t+n} - Q_{t+n-1}(S_t, A_t) \right]$$



**Figure 7.4:** Gridworld example of the speedup of policy learning due to the use of $n$-step methods. The first panel shows the path taken by an agent in a single episode, ending at a location of high reward, marked by the G. In this example the values were all initially 0, and all rewards were zero except for a positive reward at G. The arrows in the other two panels show which action values were strengthened as a result of this path by one-step and $n$-step Sarsa methods. The one-step method strengthens only the last action of the sequence of actions that led to the high reward, whereas the $n$-step method strengthens the last $n$ actions of the sequence, so that much more is learned from the one episode.

# Chapter 7: N-step Sarsa

---

**Sarsa (on-policy TD control) for estimating $Q \approx q_*$**

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$
Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+, a \in \mathcal{A}(s)$, arbitrarily except that $Q(terminal, \cdot) = 0$

Loop for each episode:
    Initialize $S$
    Choose $A$ from $S$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
    Loop for each step of episode:
        Take action $A$, observe $R$, $S'$
        Choose $A'$ from $S'$ using policy derived from $Q$ (e.g., $\varepsilon$-greedy)
        $Q(S, A) \leftarrow Q(S, A) + \alpha\big[R + \gamma Q(S', A') - Q(S, A)\big]$
        $S \leftarrow S'; A \leftarrow A';$
    until $S$ is terminal

---

**$n$-step Sarsa for estimating $Q \approx q_*$ or $q_\pi$**

Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}$
Initialize $\pi$ to be $\varepsilon$-greedy with respect to $Q$, or to a fixed given policy
Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$, a positive integer $n$
All store and access operations (for $S_t$, $A_t$, and $R_t$) can take their index mod $n + 1$

Loop for each episode:
    Initialize and store $S_0 \neq$ terminal
    Select and store an action $A_0 \sim \pi(\cdot|S_0)$
    $T \leftarrow \infty$
    Loop for $t = 0, 1, 2, \dots$ :
    |   If $t < T$, then:
    |       Take action $A_t$
    |       Observe and store the next reward as $R_{t+1}$ and the next state as $S_{t+1}$
    |       If $S_{t+1}$ is terminal, then:
    |          $T \leftarrow t + 1$
    |       else:
    |          Select and store an action $A_{t+1} \sim \pi(\cdot|S_{t+1})$
    |   $\tau \leftarrow t - n + 1$   ($\tau$ is the time whose estimate is being updated)
    |   If $\tau \geq 0$:
    |       $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
    |       If $\tau + n < T$, then $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$    $(G_{\tau:\tau+n})$
    |       $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha\left[G - Q(S_\tau, A_\tau)\right]$
    |       If $\pi$ is being learned, then ensure that $\pi(\cdot|S_\tau)$ is $\varepsilon$-greedy wrt $Q$
    Until $\tau = T - 1$

# Chapter 7: N-step off-policy Sarsa

$$V_{t+n}(S_t) \leftarrow V_{t+n-1}(S_t) + \alpha \underbrace{\rho_{t:t+n-1}}_{\text{importance sampling ratio}} [G_{t:t+n} - V_{t+n-1}(S_t)]$$

---

**Off-policy $n$-step Sarsa for estimating $Q \approx q_*$ or $q_\pi$**

Input: an arbitrary behavior policy $b$ such that $b(a|s) > 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}$
Initialize $Q(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}$
Initialize $\pi$ to be greedy with respect to $Q$, or as a fixed given policy
Algorithm parameters: step size $\alpha \in (0, 1]$, a positive integer $n$
All store and access operations (for $S_t$, $A_t$, and $R_t$) can take their index mod $n + 1$

Loop for each episode:
    Initialize and store $S_0 \neq$ terminal
    Select and store an action $A_0 \sim b(\cdot|S_0)$
    $T \leftarrow \infty$
    Loop for $t = 0, 1, 2, \dots$ :
    |   If $t < T$, then:
    |      Take action $A_t$
    |      Observe and store the next reward as $R_{t+1}$ and the next state as $S_{t+1}$
    |      If $S_{t+1}$ is terminal, then:
    |          $T \leftarrow t + 1$
    |      else:
    |          Select and store an action $A_{t+1} \sim b(\cdot|S_{t+1})$
    |   $\tau \leftarrow t - n + 1$     ($\tau$ is the time whose estimate is being updated)
    |   If $\tau \geq 0$:
    |      $\rho \leftarrow \prod_{i=\tau+1}^{\min(\tau+n, T-1)} \frac{\pi(A_i|S_i)}{b(A_i|S_i)}$               $(\rho_{\tau+1:\tau+n})$
    |      $G \leftarrow \sum_{i=\tau+1}^{\min(\tau+n, T)} \gamma^{i-\tau-1} R_i$
    |      If $\tau + n < T$, then: $G \leftarrow G + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$       $(G_{\tau:\tau+n})$
    |      $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha \rho [G - Q(S_\tau, A_\tau)]$
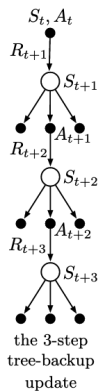    |      If $\pi$ is being learned, then ensure that $\pi(\cdot|S_\tau)$ is greedy wrt $Q$
    Until $\tau = T - 1$

# Chapter 7: Tree Backup

▶ The update rule is:

$$V_{t+n}(S_t) \leftarrow V_{t+n-1}(S_t) + \alpha\left[G_{t:t+n} - V_{t+n-1}(S_t)\right]$$



$S_t, A_t$

$R_{t+1}$

$S_{t+1}$

$R_{t+2}$ $A_{t+1}$

$S_{t+2}$

$R_{t+3}$ $A_{t+2}$

$S_{t+3}$

the 3-step
tree-backup
update

# Chapter 7: Tree Backup

- 1-step Expected Sarsa = 1-step Tree Backup

$$G_{t:t+1} = R_{t+1} + \sum_a \pi(a|S_{t+1})Q_t(S_{t+1}, a)$$

- 2-step Tree Backup

$$
\begin{aligned}
G_{t:t+2} &= R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q_{t+1}(S_{t+1}, a) \\
&\quad + \gamma\pi(A_{t+1}|S_{t+1})\left(R_{t+2} + \gamma \sum_a \pi(a|S_{t+2})Q_{t+1}(S_{t+2}, a)\right) \\
&= R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1})Q_{t+1}(S_{t+1}, a) \\
&\quad + \gamma\pi(A_{t+1}|S_{t+1})G_{t+1:t+2},
\end{aligned}
\tag{1}
$$

- n-step Tree Backup update (same as n-step Sarsa)

$$Q_{t+n}(S_t, A_t) = Q_{t+n-1}(S_t, A_t) + \alpha\left[G_{t:t+n} - Q_{t+n-1}(S_t, A_t)\right]$$

# Chapter 7: Tree Backup

**$n$-step Tree Backup for estimating $Q \approx q_*$ or $q_\pi$**

Initialize $Q(s,a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}$
Initialize $\pi$ to be greedy with respect to $Q$, or as a fixed given policy
Algorithm parameters: step size $\alpha \in (0,1]$, a positive integer $n$
All store and access operations can take their index mod $n+1$

Loop for each episode:
    Initialize and store $S_0 \neq$ terminal
    Choose an action $A_0$ arbitrarily as a function of $S_0$; Store $A_0$
    $T \leftarrow \infty$
    Loop for $t = 0, 1, 2, \dots$ :
    |  If $t < T$:
    |     Take action $A_t$; observe and store the next reward and state as $R_{t+1}, S_{t+1}$
    |     If $S_{t+1}$ is terminal:
    |        $T \leftarrow t+1$
    |     else:
    |        Choose an action $A_{t+1}$ arbitrarily as a function of $S_{t+1}$; Store $A_{t+1}$
    |  $\tau \leftarrow t+1-n$   ($\tau$ is the time whose estimate is being updated)
    |  If $\tau \geq 0$:
    |     If $t+1 \geq T$:
    |        $G \leftarrow R_T$
    |     else
    |        $G \leftarrow R_{t+1} + \gamma \sum_a \pi(a|S_{t+1})Q(S_{t+1},a)$
    |     Loop for $k = \min(t, T-1)$ down through $\tau+1$:
    |        $G \leftarrow R_k + \gamma \sum_{a \neq A_k} \pi(a|S_k)Q(S_k,a) + \gamma\pi(A_k|S_k)G$
    |     $Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha[G - Q(S_\tau, A_\tau)]$
    |     If $\pi$ is being learned, then ensure that $\pi(\cdot|S_\tau)$ is greedy wrt $Q$
    Until $\tau = T-1$

# Chapter 7: Summary